

# Seeding Influential Nodes in Non-Submodular Models of Information Diffusion

Elliot Anshelevich  
CS Department, RPI.  
eanshel@cs.rpi.edu

Ameya Hate  
CS Department, RPI.  
hatea2@rpi.edu

Malik Magdon-Ismail  
CS Department, RPI.  
magdon@cs.rpi.edu

## Abstract

We consider the model of information diffusion in social networks from [21] which incorporates trust (weighted links) between actors, and allows actors to actively participate in the spreading process, specifically through the ability to query friends for additional information. This model captures how social agents transmit and act upon information more realistically as compared to the simpler threshold and cascade models. However, it is more difficult to analyze, in particular with respect to seeding strategies. We present efficient, scalable algorithms for determining good seed sets – initial nodes to inject with the information. Our general approach is to reduce our model to a class of simpler models for which provably good sets can be constructed. By tuning this class of simpler models, we obtain a good seed set for the original more complex model. We call this the *projected greedy approach* because you ‘project’ your model onto a class of simpler models where a greedy seed set selection is near-optimal. We demonstrate the effectiveness of our seeding strategy on synthetic graphs as well as a realistic San Diego evacuation network constructed during the 2007 fires.

## 1 Introduction

Networks (social, computer, and physical) are replete with the flow of information, ideas, innovations, etc., and it is these flows which affect the way people think, act, and bind together in a society. Ideally, important messages should disseminate quickly and reach the people who need to take action, and the diffusion of malicious gossip should, if possible, be terminated. Since diffusion of both useful and malicious items is at the core of our society, it is vital to understand the mechanisms of diffusion through dynamic networks. A network can change as a result of the diffusion. For example, an evacuation message may not reach its intended audience because certain important people (critical conduits of information) left the community before the diffusion completed. In this paper, we study how to optimize a diffusion in realistic, large scale

(multi-million node) complex networks; in particular, how to select those actors to be initially seeded with information so as to maximize the ultimate number of actors receiving the information *and* acting upon that information. Of particular interest to us is the diffusion of high-value actionable information – information which is asking the user to take some action – in particular diffusion of an evacuation warning. This will be the context of our discussion, however our methods are general.

The figure on the right shows the area which was affected by fires in 2007 (shaded red), the mandatory evacuation area (yellow boundary), and regions of unwarned people (black dots). There are many unwarned people in the mandated evacuation area with at least 6 reported deaths. This leaves the question, “How can one improve the communication of such high value information?” The social network is important, and one natural avenue (given that it is not feasible to communicate to everyone) is to try to optimize with respect to a fixed budget of people who you can contact. That is, can we choose the seeds of the information diffusion to maximize the ultimate spread. This is the question we address.

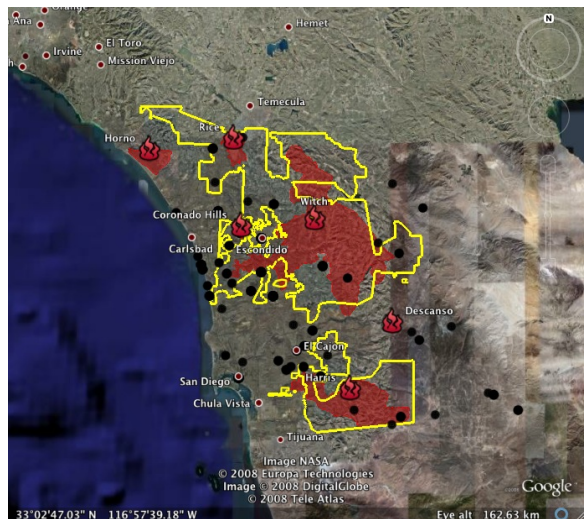


Figure 1: 2007 San Diego fires evacuation area.

The study of emergency warnings and evacuation is a good context for diffusion. There is no universally trusted news source, and even if there were, it is practically infeasible for such a news source to reach everyone. Thus, it is essential to make use of the social communication network to diffuse a warning through the network in such a way that people act. A person hearing the same information from multiple independent sources is more likely to act on it. The notion of trust, which measures the likelihood that a message from one person to another will be believed, plays an important role in such diffusions [11,12,26]. We will use the model developed in [21] to approximate the process of social information diffusion. This model captures:

- (i) the notion of trust and community structure;
- (ii) an actor’s ability to query friends for more information (people often receive warning messages from various sources such as their family and friends, the media or local authorities, through various channels such as face-to-face or telephone, and may seek confirmation and/or additional information, [51,53]);
- (iii) the existence of multiple sources of the same information, each source is trusted to a different level;
- (iv) network dynamics as a result of the diffusion (for example nodes evacuating the network);

- (v) tunable thresholds and parameters that can model different types of diffusion, from gossip in online media to evacuation warnings.

In short, we believe this model to be a reasonable model of general information diffusion on a social network. We summarize the essential details of the model in Section 2.

**The Need for Heuristics.** Given a social network and an information diffusion model, even one that is much simpler than the model hinted at above, there is no known procedure to efficiently<sup>1</sup> compute a seed set that results in a maximum number of actors taking the prescribed action. So, our general goal is to develop efficient heuristics to compute good (not necessarily optimal) seed sets on population scale social networks. This means that the heuristics should have sub-quadratic running time (ideally nearly linear running time).

## 1.1 Our Contribution

We give efficient heuristics to select a subset of the actors (the seed set) to initialize with the information with the goal of trying to maximize the final set of actors who believe and *act upon* that information. We call this a targeting algorithm. We introduce a new paradigm for developing such heuristics to construct seed sets in complex and realistic diffusion models. Our approach works as follows. First we develop an appropriate simplification of the general model (with certain tunable knobs) for which we can develop *efficient, near-optimal targeting algorithms*. In contrast to the prior work, we develop the targeting algorithms for a *different, more efficient model*, not the original general model. The key aspect of our approach is that the simpler model has tunable parameters; these tunable parameters can be adjusted so as to optimize the ultimate diffusion *for the true general model*. More specifically, there are many instances of the simpler model, each specified by a particular setting of the tunable parameters. For any such instance of the simpler model we can compute a near-optimal seed set. Since we do not know which instance of the simpler model best represents the true model, we perform a guided search through all instances of the simpler model: each instance of the simpler model gives a seed set and we pick the best one. The details are given in Section 4. Intuitively, the optimal seed sets for the simpler model give a “smart” set of seed sets for the general model, and one of these smart seeds is likely to be good.

**Experimental Testing** We demonstrate the effectiveness of our approach on several random graph models, as well as a model of the San Diego network during the 2007 fires that was created in [22]. In all cases, and for various types of diffusions (settings of parameters within the true general diffusion model)

---

<sup>1</sup>Choosing a seed set to maximize a diffusion belongs to the class of NP-hard problems, a class for which there are no known efficient procedures. A procedure is efficient if it runs in time that is polynomial in the size of the social network. For practical purposes an algorithm that takes longer than the cube of the network size is already not feasible on population scale networks. Since the only known algorithms that maximize a diffusion are exponential, such algorithms are far from feasible.

our approach gives the best seeds when compared with two benchmark algorithms: random targeting and high-degree targeting.

There is a tradeoff, however. Random targeting, though not very effective, is very easy and requires little information regarding the network other than the nodes that are present. High degree targeting requires more information, namely a measure of prominence of the nodes as well, and often, though not always, does better than random seeding. Our algorithms require the most information, namely some model for the connectivity in the network (who links to whom). Though this may not be known in practice, it can be estimated as was done in [22] for the San Diego network (based on population densities and a geometric community based random graph model that gives highest probability to links that are between geometrically close nodes within the same community and lowest probability to links that are between geometrically distant nodes within different communities).

## 1.2 Background and Related Work

In this section we give a brief summary of related work that forms a pertinent background for our research. Social networks play an important role in the spread of ideas, information, products, diseases, etc. This diffusion affects the way people think, act and make decisions in a society. Thus modeling the flow of information is an important and active research area.

Infectious diseases are similar to infectious ideas and have dynamics similar to that of evacuation warnings. Many studies use the homogeneous SIR model of Reed and Frost [40]. Both Markov chains (e.g., [10]) and differential equations [4, 41] have been used. Most of these methods make homogeneity assumptions about the underlying network, and thus do not take full advantage of the network topology.

Information diffusion is a long researched area [55], with work on online communities becoming a very active topic recently, on account of innovation diffusion, viral marketing, and computer virus spread [18, 19, 24, 25, 30, 31, 54, 56]. Two fundamental types of models of information diffusion used in the literature are cascade models and threshold models [17, 27, 28, 57]. For example, in the Independent Cascade model each node gets an opportunity to influence each of its inactive neighbors and is successful with a given probability. In the Linear Threshold model, each node has a threshold which determines the number of neighbors that need to be active for this node to become active itself. Such simple diffusion models often have mathematically convenient properties, like submodularity, which do not hold in general diffusion models. Submodularity, for instance, implies that efficient greedy algorithms can construct near-optimal seed sets [27]. Our work, in part, is to use simpler models which are submodular in order to guide the search for good seed sets in the general diffusion model.

**Warnings and Evacuation** In order to study the diffusion of warnings and evacuation messages through a population, we need to understand how various psychological, social, and economic factors drive the process. Humans have a layered warning response process that goes through multiple stages [50]. Hence, any attempt at realistically modeling this scenario must consider these factors. For this purpose we turn to work carried out in Social sciences literature.

Much research has focused on employing advanced technology in detection and prediction, or optimizing evacuation routes. Far less attention has been paid to early warning dissemination [16,33,34,49]. We conclude from warning response behavior [37] and the decision making involved [32,42,44] that it is as much social as it is about communication. The warning sequence process model [49,51,53] posits that there are six phases extending from disseminating the warning to protective action. Social science theory suggests that people’s response to warnings depends on the social, economic, demographic and physical/technological environment, [1, 14, 16, 33, 50] as well as message content [32, 33, 44, 49, 51, 53]. People may seek additional information and confirmation through observation and direct contact [37,49,51,53]. Thus, special care should be taken to incorporate these effects in any general model for diffusion of warnings. The warning time distribution has two components: the official broadcast component [44,47,52] and the information contagion (diffusion) component [14,47,51,53]. There has been considerable research in estimating warning time distributions by modeling the warning network, however the contagion component is little understood [1,35,44,47], and needs to be better understood for protective agencies to fully estimate warning times [33,34]. Any tool that wants to analyze warning time distribution must incorporate the effects of inhomogeneities such as social groups, ethnic groups, and in general differential trust in and access to warning systems [35,45,46,51,53]. The model of diffusion we use captures trust as well as human reactivity and was built around these social considerations. Our contribution is *not* the model, which we take as given from prior work [21]. Our contribution is to take this social-science based model of diffusion and construct good targeting algorithms for it.

**Agent based models** In scenarios where individuals are influenced by their social environment, agent based modeling is often used [13,18]. In such cases, decision making entities are represented by agents whose behavior can be based on a set of rules. Such a model allows the agent to have intelligence and memory and also exhibit complex behaviors like learning and adapting [8,38]. In the context of diffusion, agent based modeling has been used, for example, to study epidemic spread [6,43]; to model information diffusion in virtual marketplace [39]; and to simulate technological diffusion [3,36] and environmental innovations [7,48]. Our research is not in the development of such models but in the exploitation of such models in determining optimal seed sets.

### 1.3 Outline of the Paper

In the next section, we summarize the diffusion model from [21] and give some analysis of its complexity in Section 3. We present our projected greedy heuristic in Section 4. In Section 5 we describe the experimental design we used in extensively evaluating our projected greedy heuristic on both synthetic networks as well as the San Diego network. We end in Section 6 where we give an overview of the results together with a discussion.

## 2 The Diffusion Model

We summarize the diffusion model that was developed in [21]. This model was used to study propagation of evacuation news through a population in [21]. We use it as our model of diffusion because it is general and contains many other simpler models as sub-cases. The model takes into consideration the fact that agents may act on information and leave the network; this means diffusion occurs on a network that is dynamic. The model contains the concepts found in the widely studied SIR model of epidemiology, as well as standard threshold and cascade models together with actions a social agent may perform, such as “querying for more information” as well as notions of trust.

We have an initial graph  $G = (V, E)$  where  $V$  are individuals in the population and the edges  $(u_i, u_j) \in E$  represent social connections between individuals  $u_i$  and  $u_j$ . There are  $K$  sources with *information values*  $I_1, I_2, \dots, I_K$ . Also, a source  $k \in \{1, \dots, K\}$  can seed  $B_k \geq 0$  nodes. We will use  $k$  to index sources and  $i, j$  to index nodes.

Associated with edge  $(u_i, u_j) \in E$  is a trust value  $0 \leq \alpha(u_i, u_j) \leq 1$ . It represents the amount of trust node  $u_j$  has on information provided by  $u_i$ . The graph may be directed with different trust values on edges  $(u_i, u_j)$  and  $(u_j, u_i)$  representing asymmetrical trust. Each node  $u \in V$  also has a similar trust value for each source  $k \in \{1, 2, \dots, K\}$ :  $\alpha(u, k)$ .

Each node  $u$  possesses an *information-value set*  $S(u) = \{(s_1, v_1), \dots, (s_K, v_K)\}$  for the  $K$  information sources. It is made up of pairs  $(s_k, v_k)$ , where  $v_k$  is the information node  $u$  has from source  $s_k$ . Based on its information value set, a node calculates its *information value* as follows:

$$I(u) = \lambda_d \sum_{k=1}^K v_k + (1 - \lambda_d) \cdot \max_{k=1, \dots, K} v_k, \quad (1)$$

where  $0 \leq \lambda_d \leq 1$  is a model parameter. The information value is a convex combination of the total information the node has from all the sources and the maximum value the node has from any one source; as such,  $I(u)$  is at most this total information, and at least the maximum (a node’s information is at least the information of its best source and at most the sum of all its sources). The two extremes  $\lambda_d = 0$  and  $\lambda_d = 1$  correspond to two different extremes of a diffusion. When  $\lambda_d = 0$ , a node ignores all information

but its highest information value source; this corresponds to a very conservative choice and would be more appropriate for something like a warning. When  $\lambda_d = 1$ , the information value is the sum, i.e., the different sources reinforce each other. This is a more aggressive diffusion, appropriate for something like gossip or a rumor.

Associated with each node  $u$  are two thresholds, a lower threshold  $t_l(u)$  and an upper threshold  $t_h(u)$  such that  $0 \leq t_l(u) \leq t_h(u)$ . The thresholds determine how a node acts given its information set. Based on its information value, a node  $u$  lies in one of the following states:

1. **Disbelieved (low information value, less than the lower threshold):** If  $I(u) < t_l(u)$  then the node does not believe it has any meaningful information. In this state, node  $u$  does not take any action except for incorporating information-value sets that it (node  $u$ ) receives from its neighbors.
2. **Undecided (intermediate information value):** If  $t_l(u) \leq I(u) < t_h(u)$  then node  $u$  believes it has some information but is uncertain about acting (which in our case is to evacuate). In this state, node  $u$  will *query* the information-value sets of its neighbors and incorporate any new information into its own information set.
3. **Believed (high information value, above the upper threshold):** If  $I(u) \geq t_h(u)$  then the node has enough information to accept the information as correct and act upon it. In this state, node  $u$  actively propagates its information-value set to its neighbors for  $\tau$  time-steps before evacuating. Here  $\tau > 0$  is a model parameter that determines the time it takes a believed node to evacuate the network (in the case of a warning message). If  $\tau = \infty$ , the node never leaves the network and propagates its information forever.
4. **Evacuated:**  $\tau$  time-steps after a node enters the Believed state, it evacuates resulting in its disconnection from the network. Hence the graph changes and a new graph  $G'$  is obtained. The new graph  $G'$  is the induced subgraph on the vertices  $V \setminus u$ ; that is  $G' = (V', E')$ , where  $V' = V \setminus u$  and  $E' = E \setminus \{(u, u_i) \mid u_i \in V', (u, u_i) \in E\}$ .

## 2.1 The Diffusion Process

Initially, all nodes have zero information, so  $v_k = 0 \ \forall k \in \{1, \dots, K\}$  and all nodes are in the disbelieved state. Each node's information set is zero,

$$S(u) = \{(s_1, 0), \dots, (s_K, 0)\}.$$

The diffusion process begins when some nodes are *seeded* with information from sources according to a seeding strategy. Seeding entails transfer of information from sources to the information-value set of the

selected seed nodes. This transfer of information gets attenuated by the trust between the node and the source. So if source  $k$  seeds node  $u$  then the information transfer to node  $u$  has value  $\alpha(u, k) \cdot I_k$  and the pair  $(s_k, 0)$  in the node's information-value set gets updated to  $(s_k, \alpha(u, k) \cdot I_k)$ . Each source seeds some subset of the nodes in this way, and multiple sources can seed the same node. After seeding, each node then computes its information value using `eq:lu`, ascertains its state, and performs the required action at the next step. At every consecutive step nodes take action depending on their state (which may involve broadcasting its information-value set and/or receiving information-value sets from its neighbors); if any new information is received, this information is merged into its current information-value set. In this way a node's information-value set evolves as the diffusion process proceeds.

We now describe the process of propagating and updating information sets. When node  $u_i$  propagates its information-value set to neighbor  $u_j$ , the value of information gets similarly attenuated by a factor of  $\alpha(u_i, u_j)$ . If

$$S(u_i) = \{(s_1, v_1), \dots, (s_K, v_K)\},$$

then the information-value set that is received by  $u_j$  is

$$\alpha(u_i, u_j) \cdot S(u_i) = \{(s_1, \alpha(u_i, u_j) \cdot v_1), \dots, (s_K, \alpha(u_i, u_j) \cdot v_K)\}.$$

Consider a source  $k$ , and consider node  $u_i$ . The node  $u_i$  may receive information value originating from source  $k$  either directly from source  $k$  or indirectly from one of its neighbors. Suppose that node  $u_i$  has  $\delta$  neighbors. In principle  $u_i$  can receive (either in the current step or in some prior step) information-value  $v_k^0$  (directly from the source) and information values  $v_k^1, v_k^2, \dots, v_k^\delta$  from each of its neighbors (each of these information values will be the attenuated value that was propagated to  $u_i$ ; some or all of these values could be zero). In order to determine its information value for this source  $k$ , the node needs to *fuse* all these information values into a single value as follows

$$v_k = \lambda_s \cdot \sum_{j=0}^{\delta} v_k^j + (1 - \lambda_s) \cdot \max_{j=0, \dots, \delta} v_k^j. \quad (2)$$

Again  $\lambda_s$  is a model parameter such that  $0 \leq \lambda_s \leq 1$ . This fusion of information happens for every source  $k$  at every node  $u_i$ . As with  $\lambda_d$ ,  $\lambda_s$  impacts how aggressive a diffusion is. At the two extremes:  $\lambda_s = 0$ , a node takes the maximum value it hears about the information from a source (the conservative case); and, when  $\lambda_s = 1$  a node takes all the information it hears about a source and adds (the aggressive diffusion). Thus, at each consecutive time step, every node updates its information-value set based on the new information that was propagated to it. A node calculates its new information value based on this updated set and this updated information value will be used to determine the node's state and possible action. The diffusion process continues, i.e. information continues to propagate as described above, until either all nodes evacuate or there is no change in the information value sets of the nodes.



**Instance  $\mathbb{G}$** 

Graph  $G = (V, E)$  specifying the network for the diffusion.

Source information values and budgets,  $(I_1, B_1), \dots, (I_K, B_K)$

Trust values  $\alpha(u_i, u_j)$  for all edges  $(u_i, u_j) \in E$ .

Trust values  $\alpha(k, u_j)$  between each source  $k$  and each node  $u_j \in V$ .

Diffusion parameters  $\lambda_d, \lambda_s, \tau$

Lower and upper thresholds  $t_l(u), t_h(u)$  for each node  $u \in V$ .

**Desired output:** seed sets  $\psi_1, \dots, \psi_K$  for each source  $k = 1, \dots, K$ .

Figure 2: An instance of the general diffusion model.

The graph  $G$  is typically given and chosen to model some social network, as for example the San Diego network during the 2007 fires. The model parameters,  $(\lambda_s, \lambda_d)$ , the evacuation time  $\tau$ , the threshold parameters  $t_l(u)$ ,  $t_h(u)$  at each node  $u$  and the edge trust values  $\alpha(u_i, u_j)$  can be chosen to model various types of diffusion settings. For example in a social network with strong communities, the trust values of edges within a community would be high (close to 1) and the edges between communities would typically be low (close to 0). In a network that has been “primed”, for example a community that has recently experienced a tsunami, the community may be in a “panic” state, which could be modeled by very low upper thresholds  $t_h(u)$ . One would also set the thresholds to be low for low risk diffusions like gossip. However for actions that incur significant cost, like a node evacuating, the lower threshold may be low but the upper threshold would typically be high. Lastly the  $\lambda$ -parameters could be chosen to model fast or slow diffusing information. In all cases, however, all these parameters are exogenously specified. An instance of the general diffusion model  $\mathbb{G}$  is summarized in Figure 2. Our goal is to optimally seed the diffusion given the graph and the parameter settings.

## 2.2 Seeding the Diffusion

The goal of this study is to find a seeding strategy that maximizes the number of nodes that end up in their Believed state, given the diffusion setting as described in the previous section, and a seeding budget, as we describe here.

Each source  $k$  has a budget  $B_k$  of nodes which it can seed. Our task is to determine the seed sets  $\psi_k \subset V$ , where  $|\psi_k| \leq B_k$ . The seed set  $\psi_k$  specifies which nodes are seeded by source  $k$ . The objective is to maximize the number of nodes that become believers. Thus, there is some function  $\Gamma$ , which we call the *coverage function*, that computes the number of nodes which become believers, given the diffusion setting and the seed sets  $\psi_1, \dots, \psi_K$ . More generally, if there is some randomization in the communication (propagation

of information-value sets)<sup>2</sup>, then the coverage function would compute the expected number of believers. Thus, in general, the coverage function maps  $(\psi_1, \dots, \psi_K)$  to  $\mathbb{R}_{\geq 0}$ .

One of the challenges is to efficiently compute the coverage function  $\Gamma$  for a given input seeding  $(\psi_1, \dots, \psi_K)$ . The other, which is our main goal, is to find the seeding which maximizes  $\Gamma$ .

### 3 Analysis

It is useful to have a theoretical understanding of the diffusion process in order to identify where the potential difficulties lie when choosing an optimal seed set. In fact, as we will soon see, the general model described in the previous section is extremely difficult for theoretical analysis. As a result, we will consider a simplified instance of the diffusion model in our theoretical analysis, and use the insight from this analysis to develop a seeding algorithm for the general model. When testing the algorithm, however, we will use instances of the general diffusion process (see Section 2).

#### 3.1 Monotonicity and Submodularity

The complexity of seed selection for the general diffusion model results from the behavior of the coverage function  $\Gamma$  when you perturb the seed set. If  $\Gamma$  behaves well when you perturb the seed set in certain ways, then simple, iterative greedy algorithms are effective at selecting near-optimal seed sets.

**Monotonicity.** The coverage function  $\Gamma$  is monotone if adding nodes to a seed set can only increase the coverage function's value. Mathematically, if  $\psi \subseteq \psi'$ , then

$$\Gamma(\psi, \psi_{-i}) \leq \Gamma(\psi', \psi_{-i}). \quad (3)$$

(We use the notation  $\Gamma(\cdot, \psi_{-i})$  to denote the coverage function as a function of its  $i$ th set argument, keeping all the other sets fixed.) It is intuitive that if you seed more nodes in an evacuation scenario, then more people should evacuate.

**Submodularity.** Submodularity captures the intuitive notion of diminishing returns. As you seed more and more, the additional benefit you get in terms of the increased coverage is decreasing. One of the equivalent mathematical definitions of submodularity is as follows.  $\Gamma$  is submodular, if for all sets  $\psi \subseteq \psi'$  and any  $A \subseteq V$ ,

$$\Gamma(\psi \cup A, \psi_{-i}) - \Gamma(\psi, \psi_{-i}) \geq \Gamma(\psi' \cup A, \psi_{-i}) - \Gamma(\psi', \psi_{-i}), \quad (4)$$

---

<sup>2</sup>For example if, when a node  $u_i$  propagates an information-value set to  $u_j$ , it is received at the other end with some probability  $p(i, j)$  depending on the communication infrastructure, then the number of nodes which are ultimate believers is a random variable.

that is the increase in  $\Gamma$  from  $\psi' \rightarrow \psi' \cup A$  is not more than its increase from  $\psi \rightarrow \psi \cup A$ . Submodularity is essentially the set-function version of concavity.

A coverage function that is submodular has the nice property that if you use a simple greedy strategy to select a seed set, then the resulting seed set has a coverage which is within a small constant factor of optimal. Unfortunately, as we are about to demonstrate, the coverage function for our general diffusion model is not submodular. It is not even monotone. The next two examples illustrate why. In both examples we set  $\lambda_d = \lambda_s = 0$ , so we are using the max function for both the computation of the information value and for the fusion of information received from neighboring nodes on the same source; we sometimes call this the *max-max* model. We have only one source ( $K = 1$ ) and its information value is  $I_1 = 1$ .

**Example 1: non-monotonicity:** The diffusion setting is shown in Figure 3. There is one low trust edge (indicated in red) and a single “high strung” node  $c$  which has very low thresholds (also indicated in red); we will soon see why we call this node high strung.

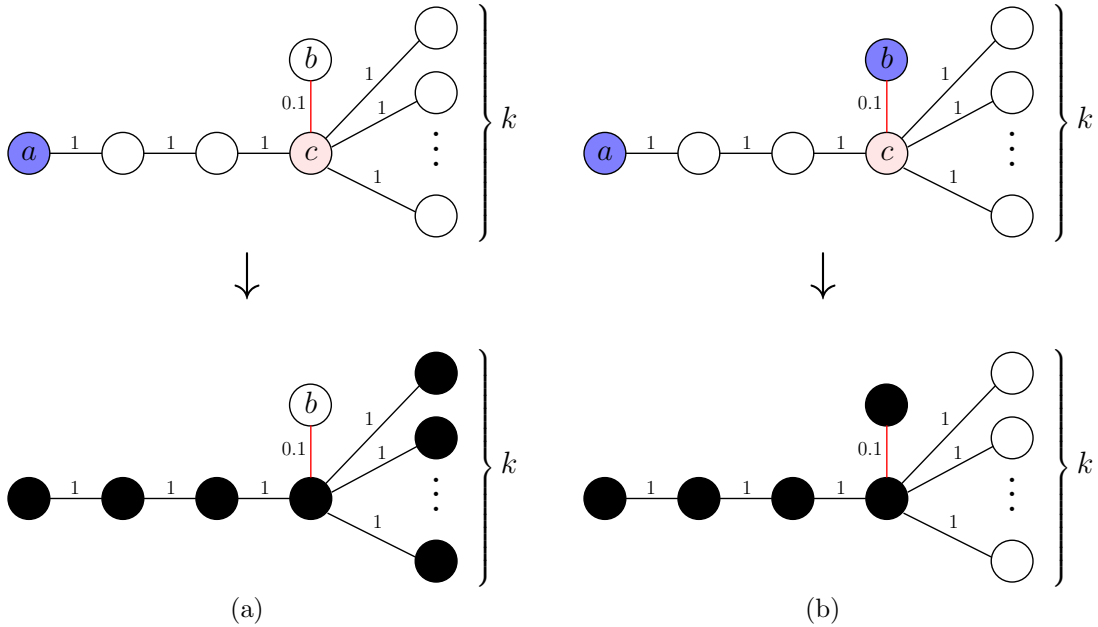


Figure 3: Non-monotonic diffusion setting. The blue nodes are the initial seed set. The black nodes are the evacuated nodes in the final state. All thresholds are 0.5 except for the high strung node  $c$  with threshold 0.1. The number of steps to evacuation is  $\tau = 1$ .

Consider the coverage function for this diffusion setting and two different seed sets  $\psi = \{a\}$  (case (a) in Figure 3) and  $\psi' = \{a, b\}$  (case (b) in the Figure 3) – in Figure 3) the seed sets are shaded blue for both cases. In case (a), the information value of 1 will propagate unattenuated to node  $c$  and beyond to the

additional  $k$  nodes. All these nodes will evacuate. Only an information value of 0.1 propagates to node  $b$  which will not evacuate. So,  $\Gamma(\psi) = k + 4$ . Now consider case (b) with seed set  $\psi'$ . Initially node  $c$  will receive attenuated information of value 0.1 from node  $b$ . Since node  $c$  has very low thresholds, in a sense it panics, and leaves the network almost immediately ( $\tau = 1$ ), after propagating its information value of 0.1. Unfortunately, however, this value of 0.1 is not high enough to evacuate any of the  $k$  peripheral nodes, since their thresholds are high. The nodes between  $a$  and  $c$  are fine, though, because they will eventually receive an information value of 1 propagated from  $a$ . The important thing is that when  $c$  leaves, it cuts off the  $k$  peripheral nodes from  $a$  which results in only 5 nodes evacuating. Hence,  $\Gamma(\psi') = 5$ .

The high strung node  $c$  will leave if it gets only a small amount of information. If, as in this case, the high strung node is crucial (is a “bridge” node in the network), then when the high strung node leaves, it disconnects potentially large parts of the network from the information before the higher value information has a chance to flow into those parts of the network. By increasing the seed set, one might increase the chances that such a high strung node gets low-value information too early in the diffusion, and the ensuing early evacuation of this high strung node is what leads to the non-monotonicity of the coverage function. In fact, if nodes do not leave the network until all information propagation has occurred, then monotonicity is guaranteed. As the next example will show, however, even if nodes do not leave the network, the coverage function is still not submodular.

**Example 2: non-submodularity.** The diffusion setting is shown in Figure 4. We set the evacuation time  $\tau = 10$  (i.e., large enough so that nodes only start leaving after all information has propagated).

Now consider the seed sets  $\psi = \emptyset$ ,  $\psi' = \{a\}$ . Clearly  $\Gamma(\psi) = 0$ . As for  $\psi'$ , if you seed  $a$ , an attenuated information of 0.9 reaches  $d$  which is not enough to breach the threshold of 0.91. Hence only  $a$  evacuates and so  $\Gamma(\psi') = 1$ . We now consider  $\psi \cup \{b\}$  and  $\psi' \cup \{b\}$ . We will see that

$$\Gamma(\psi \cup \{b\}) - \Gamma(\psi) < \Gamma(\psi' \cup \{b\}) - \Gamma(\psi'), \quad (5)$$

which contradicts submodularity, since  $\psi \subseteq \psi'$ . First consider  $\psi \cup \{b\} = \{b\}$  which is case (a) in Figure 4. An information value of 0.1 propagates to  $c$  which will therefore enter the undecided state since  $t_l(c)$  is low. However, the neighbors cannot provide any new information so nothing further will happen. Therefore, only  $b$  evacuates.

Now consider  $\psi' \cup \{b\} = \{a, b\}$  which is case (b) in Figure 4. As before, 0.9 propagates to  $c$  and 0.9 to  $d$ . But, because  $c$  is in the undecided state, it will *query*  $d$  and receive information of 0.9 which is enough to push  $c$  to become a believer. Now,  $c$  will propagate its information value of 0.9 to  $e$  and beyond. Node  $d$ , however, will always remain below its threshold of 0.91. Hence  $\Gamma(\psi' \cup \{b\}) = k + 4$ . So

$$\Gamma(\psi \cup \{b\}) - \Gamma(\psi) = 1 < k + 3 = \Gamma(\psi' \cup \{b\}) - \Gamma(\psi').$$

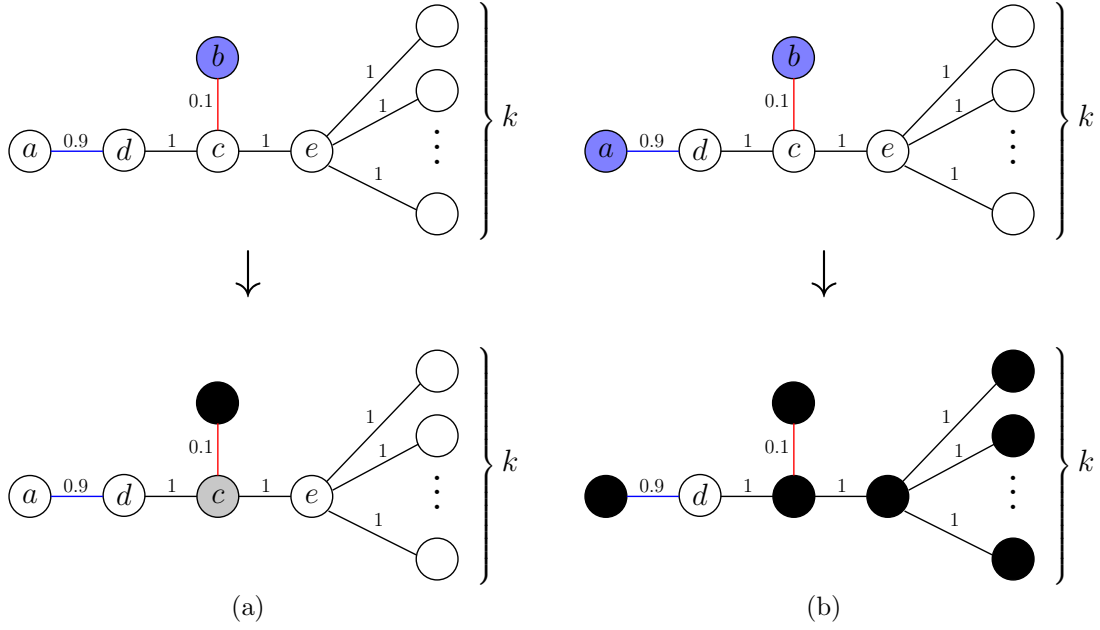


Figure 4: Non-submodular diffusion setting. The blue nodes are the initial seed set. In the final state, the black nodes are the evacuated and the light gray are undecided. All thresholds are 0.5 except for nodes  $c$  and  $d$ : for node  $d$ ,  $t_l(d) = t_h(d) = 0.91$  and for node  $c$ ,  $t_l(c) = 0.1$ ,  $t_h(c) = 0.9$ . The number of steps to evacuation is  $\tau = 10$ .

The fact that node  $c$  could query for information that  $d$  did not transmit is what resulted in the cascade effect which converted node  $e$  and beyond to the Believed state. Though it is natural for humans in a social diffusion process to query, it is precisely this ability to query that leads to non-submodularity and the resulting complexity of the process.

### 3.2 A Simplified Model

Taking a cue from the examples above, we define a simplified model that is a strict subcase of the general diffusion model defined earlier. Note that while doing experimental simulations we will use the general diffusion model. We use the simplified model only to carry out theoretical analysis. Insights obtained for this simplified model help us develop heuristics that we then apply to the general case.

In this simplified model  $\lambda_d = \lambda_s = 0$ . So, the information value at a node is the maximum value among source-value pairs, i.e.,  $I(u) = \max_{k=1,\dots,K} v_k$ ; similarly, when fusing information arriving via different neighbors regarding the same source we take the maximum; we call this the *max-max* model. We also assume all nodes have a single threshold, so  $t_l(u) = t_h(u) = t(u)$ . This eliminates the Undecided state which, as we saw, leads to non-submodularity in the coverage function. To ensure monotonicity, we choose the evacuation

time  $\tau$  to be sufficiently large (e.g.,  $\tau \gg |E|$ ). This means that nodes do not leave the network until all information propagation has occurred. For this simplified *max-max* model, we prove a simple lemma that characterizes the nature of the coverage function. We are interested in

$$\Gamma(\psi_1, \dots, \psi_K).$$

Define the *singleton* coverage (set) function  $\gamma(u, k)$  to be the *set* of nodes ultimately converted to the Believer state when only one node  $u$  is seeded by a source  $k$  with information value  $I_k$ . Thus,  $\Gamma(\psi_1, \dots, \psi_K)$  equals  $|\gamma(u, k)|$  when  $\psi_i = \emptyset$  for all  $i \neq k$ , and  $\psi_k = \{u\}$ . The next lemma characterizes the form of  $\Gamma(\psi_1, \dots, \psi_K)$ . It basically says that to obtain the set of nodes converted to the Believer state, it suffices to consider each source and let it just seed one of the nodes in its seed set, with no other source seeding any nodes. Some set of nodes is converted to Believer in this process. We consider, in this way, all the (source, seed) pairs in turn, computing the converted set when just this source seeds just this one node. By taking the union of all these converted sets, we get the final set of nodes converted when all the sources simultaneously seed all the nodes in their respective seed sets.

**Lemma 1** *The set of nodes converted to Believer with seed sets  $\psi_1, \dots, \psi_K$  is*

$$\bigcup_{k=1}^K \bigcup_{x \in \psi_k} \gamma(x, k).$$

**Proof.** The lemma follows because we are using the *max-max* model. First, suppose that  $u \in \gamma(x, k)$ . Nodes receive information in time steps. Let  $I_m(u)$  be the information value at node  $u$  after time step  $m$  of information propagation using a single seed  $(x, k)$ . Similarly let  $I'_m(u)$  be the information value at  $u$  after time step  $m$  of information propagation using all sources and their respective seeds. We claim that  $I'_m(u) \geq I_m(u)$  for all  $m \geq 0$  and all  $u$ . Indeed, suppose to the contrary that this does not hold for some  $m, u$ ; in which case, there is an earliest time step  $m$  for which it does not hold, i.e.  $I'_m(u) < I_m(u)$  and for all  $\ell < m$  and all  $v$ ,  $I'_\ell(v) \geq I_\ell(v)$ . It means that  $I_m(u)$  is different from  $I_{m-1}(u)$ . Since there is no querying, and since we are using the *max-max* model, it means that some node  $v$  *propagated* information to node  $u$  in the previous time step that resulted in information value  $I_m(u) = \alpha(v, u)I_{m-1}(v)$ . But by assumption,  $I'_{m-1}(v) \geq I_{m-1}(v)$  and since we are using the *max-max* model,  $I'_m(u) \geq \alpha(v, u)I'_{m-1}(v) \geq \alpha(v, u)I_{m-1}(v) = I_m(u)$ , a contradiction. Thus, if a node is converted in any set  $\gamma(x, k)$ , then it is also converted in the joint information propagation using all the seed sets simultaneously.

Let  $I_m^{\max}(u)$  denote the maximum value of  $I_m(u)$  over all singleton (seed, source) pairs  $(x, k)$  where  $x \in \psi_k$  and  $k \in \{1, \dots, K\}$ . We now claim that  $I'_m(u) \leq I_m^{\max}(u)$ . This means that if a node gets converted in the joint information propagation, it gets converted in at least one of the singleton information propagations. Indeed, suppose to the contrary that for some earliest  $m$ ,  $I'_m(u) > I_m^{\max}(u)$ . Again this means that in the joint propagation,  $I'_m(u)$  changed at step  $m$ , and so information was propagated from some node  $v$  to  $u$

with the result that  $I'_m(u) = \alpha(v, u)I'_{m-1}(v)$  (*max-max* model). We also know that there is some singleton propagation with  $I_{m-1}(v) \geq I'_{m-1}(v)$ , by assumption, since  $m$  is the earliest time step when this fails. In this singleton propagation, it must be that  $I_m(u) \geq \alpha(v, u)I_{m-1}(v) \geq \alpha(v, u)I'_{m-1}(v) = I'_m(u)$  (*max-max* model), a contradiction. ■

An immediate consequence of Lemma 1 is that if all the source values are the same, so  $I_k = I$  for  $k \in \{1, \dots, K\}$ , and every node trusts all the sources equally, so  $\alpha(k, u)$  is independent of  $k$ , then the converted set is

$$\bigcup_{k=1}^K \bigcup_{u \in \psi_k} \gamma(u, 1).$$

This converted set is exactly what it would be if there was just one source with value  $I$  seeding  $\psi = \bigcup_{k=1}^K \psi_k$ .

**Lemma 2** *If all sources are identical and every node trusts all sources equally, then  $K$  sources with information value  $I$  seeding sets  $\psi_1, \dots, \psi_K$  results in the same converted set as one source with information value  $I$  seeding the set  $\psi = \bigcup_{k=1}^K \psi_k$ .*

It is therefore immediate that finding the optimal seed sets for  $K$  identical sources with budgets  $B_1, \dots, B_K$  is equivalent to finding the optimal single seed set for just a single one of these sources with budget  $\sum_{k=1}^K B_k$ . Thus, the  $K$  identical sources problem reduces to the single source problem.

A direct application of Lemma 1 gives the coverage function.

**Lemma 3**

$$\Gamma(\psi_1, \dots, \psi_K) = \left| \bigcup_{k=1}^K \bigcup_{u \in \psi_k} \gamma(u, k) \right|.$$

It follows from Lemma 3 that  $\Gamma(\psi_1, \dots, \psi_K)$  is monotone. It is also easy to see that functions of this form (that are the size of the union taken nodewise of a set function defined on a node) are submodular. We therefore have the following theorem.

**Theorem 1** *For the max-max model, there is a greedy deterministic algorithm which computes seed sets  $\psi_1, \dots, \psi_K$  of sizes  $B_1, \dots, B_K$  for which*

$$\Gamma(\psi_1, \dots, \psi_K) \geq \frac{1}{2} \cdot \Gamma(\psi'_1, \dots, \psi'_K) \quad \forall \psi'_k \text{ with } |\psi'_k| \leq B_k.$$

*Moreover, when the individual budget constraints  $|\psi_k| \leq B_k$  are replaced by a total budget constraint  $|\psi_1| + \dots + |\psi_K| \leq B$ , then the deterministic greedy algorithm yields seed sets  $\psi_1, \dots, \psi_K$  of total size equal to the budget  $B$  such that*

$$\Gamma(\psi_1, \dots, \psi_K) \geq \left(1 - \frac{1}{e}\right) \Gamma(\psi'_1, \dots, \psi'_K) \quad \forall \psi'_k \text{ with } |\psi'_1| + \dots + |\psi'_K| \leq B.$$

- 1: **Input:** Instance of the *max-max* model with  $K$  sources and budgets  $B_1, \dots, B_K$ .
- 2: Initiate converted nodes  $C = \emptyset$  and selected nodes  $\psi_k = \emptyset$ .
- 3: Compute  $C_{u,k} = \gamma(u, k)$  for all nodes  $u \in V$  and sources  $k = 1, \dots, K$ ;
- 4: **while**  $|\psi_k| < B_k$  for any  $k$  and  $|C| < |V|$  **do**
- 5:   Choose  $(u^*, k^*)$  with  $|\psi_{k^*}| < B_{k^*}$  that maximizes  $|C_{u,k} \setminus C|$  over  $(u, k)$ ;
- 6:   Update  $\psi_k \leftarrow \psi_k \cup u^*$  and  $C = C \cup C_{u^*, k^*}$ ;
- 7: Output sets  $\psi_1, \dots, \psi_K$ .

**Algorithm 1:** Greedy algorithm on instance of simplified model

**Proof.** The theorem follows from the monotonicity and submodularity of  $\Gamma$  implied by Lemma 3 and the greedy deterministic algorithm for maximizing a monotone submodular function ([15]). The case with a total budget simply asks to find a set of (node,source) pairs of size at most  $B$  maximizing a submodular monotone set function; for this case the greedy algorithm gives a  $1 - \frac{1}{e}$  approximation. The case with multiple sources asks to maximize a submodular monotone set function with respect to partition matroid constraints; for this, the greedy algorithm from [15] gives a  $\frac{1}{2}$ -approximation. ■

**Remark.** When  $K = 1$  (single source), the total budget constraint and the individual budget constraints are equivalent and we get a  $(1 - \frac{1}{e})$ -approximation.

**Remark.** Randomized algorithms guaranteeing a  $(1 - \frac{1}{e})$ -approximation exist for the case with multiple sources as well as for a single source [9]. We choose to focus on the greedy algorithms due to their simplicity and relative efficiency, as our goal is to find good seed sets for extremely large networks.

The greedy algorithm implied by Theorem 1 is an intuitive algorithm that at each greedy step selects the node to add to the seed set that gives the largest increase in the size of the converted set. This algorithm is summarized in Algorithm 1 for the case with individual budget constraints. When there is a total budget constraint, the algorithm is exactly the same except that the **while** condition in step 4 checks that the total budget constraint is not exceeded. Our projected greedy heuristic for the general model is based on this greedy algorithm. Implementation of this algorithm is discussed in Section 4.



## 4 The Projected Greedy Heuristic

We now describe the Projected Greedy heuristic which takes as input an instance  $\mathbb{G}$  of the general diffusion model (see Figure 2) with  $K$  sources and produces as output a seed set  $\{\psi_1, \psi_2, \dots, \psi_K\}$ . Given the seed sets  $\{\psi_1, \psi_2, \dots, \psi_K\}$ , it is possible to compute the coverage set by simulating the information set updates in the network; simulating a single time-step takes  $O(|E|)$ . We will denote the time it takes to compute the coverage set in the general model by  $T_{\mathbb{G}} = O(a|E|)$ , where  $a$  is the number of timesteps that it takes for the diffusion process to converge. While theoretically  $a$  can be large for the general model of diffusion we consider, practically we have observed in our experiments that 30-50 timesteps is sufficient. Thus typically  $T_{\mathbb{G}}$  is on the order of the number of links in the network, which for sparse social network graphs is  $O(|V|)$ . This is useful, because all our algorithms need to be able to evaluate a seeding in order to improve it. First, for comparison, we describe two natural algorithms because our projected greedy algorithm has aspects of both.

### 4.1 Brute Force

The brute force algorithm is extremely simple: try every possible distinct seeding, compute  $\Gamma$  for each seeding and select the seeding which maximizes the coverage. The running time of this naive brute force approach is  $O\left(T_{\mathbb{G}} \prod_{k=1}^K \binom{|V|}{B_k}\right)$  because there are  $\prod_{k=1}^K \binom{|V|}{B_k}$  possible seed sets. Clearly, this algorithm is not feasible, even for  $B_k = 1$  (there are just too many possible seedings to test). However, if we could somehow intelligently prune this collection of seedings to a small number, then it would become viable. This is one aspect of the projected greedy approach: obtain plausible candidate seedings and pick one of them using the brute force approach.

The way we will obtain these plausible candidate screenings is using a greedy approach. The natural greedy algorithm is what we describe next, which we call *Actual Greedy*. Actual greedy deterministically produces a single seed set.

### 4.2 The Actual Greedy Approach

The natural greedy seeding strategy is also simple. For every pair  $(u, k)$  where  $k$  is a source and  $u \in V$  a node, we consider adding it to the seed set (as long as we do not violate the budget constraints  $|\psi_k| \leq B_k$ ). We add the pair  $(u^*, k^*)$  which results in the largest increase in the coverage  $\Gamma(\psi_1, \dots, \psi_K)$ . Starting with an empty seeding, in this way we build up the seeding to the final seeding one pair at a time. We call this strategy *Actual Greedy*. While *Actual Greedy* is a plausible heuristic, because a general instance of the diffusion model is non-monotone and non-submodular, there is no performance guarantee as compared with using the greedy algorithm for the simplified model (Section 3.2).

**Running Time of *Actual Greedy*.** Each step of *Actual Greedy* adds a node to the seeding. If the total budget is  $B = \sum_{k=1}^K B_k$ , then there are  $B$  steps in the algorithm. At each step, we need to consider  $O(K|V|)$  pairs to determine which one is the best to add; and the time to test the seeding with that pair added is  $T_{\mathbb{G}}$ . So the total running time is  $O(B \cdot K \cdot |V| \cdot T_{\mathbb{G}})$ . For sparse graphs,  $|E| = O(|V|)$ , and when the total size of the seed set is a fraction of the graph ( $B = \Theta(|V|)$ ), this means that the running time is typically  $O(K|V|^3)$ .

The general diffusion model for arbitrary settings of the parameters does not have any nice properties that can be exploited to improve this running time, and so for large graphs, with millions of nodes, this cubic running time is practically infeasible.

### 4.3 Projected Greedy

Given a general instance of the diffusion model  $\mathbb{G}$ , the idea behind the *Projected Greedy* algorithm is to construct an instance of the simplified model  $\mathbb{S}$  (see Section 3.2) that closely approximates  $\mathbb{G}$ . In a matter of speaking, we are “projecting”  $\mathbb{G}$  down to the simpler instance  $\mathbb{S}$  that most “closely” approximates  $\mathbb{G}$ . For the simpler instance  $\mathbb{S}$ , we can leverage Theorem 1 and use the greedy algorithm to obtain a constant factor approximation to the optimal seeding in  $\mathbb{S}$ . If the instance  $\mathbb{S}$  is a close approximation to the instance  $\mathbb{G}$  then the near optimal seeding for  $\mathbb{S}$  will also be a good seeding for  $\mathbb{G}$ .

There are two advantages of *Projected Greedy* over *Actual Greedy*: (1) Since the simplified model is monotone and submodular, we can use a more efficient algorithm than running the  $O(K|V|^3)$  greedy algorithm on the general instance  $\mathbb{G}$ ; (2) Neither *Actual Greedy* nor *Projected Greedy* give any performance guarantee for the quality of the seeding. However, we do have some flexibility in the choice of the simplified instance  $\mathbb{S}$ . So, by exploring a variety of plausible simplified instances, we can generate several seedings (say  $c$  of them), and we can choose one of these  $c$  as in the Brute Force approach, via simulation. One can view this as a more intelligent sampling of the possible seedings to test in the Brute Force approach, where the choice of seedings is guided by the various choices for the simplified instances  $\mathbb{S}$ .

Although we cannot give a guarantee on the quality of the seeding produced by our heuristic, our experimental evaluation demonstrates that *Projected Greedy* performs significantly better than widely used simple seeding strategies.

### 4.4 Creating an Instance of the Simplified Model

Given an instance  $\mathbb{G}$  of the general diffusion model (see Figure 2), we construct an instance  $\mathbb{S}$  of the simplified model by setting  $\lambda_d = \lambda_s = 0$ : an instance of the *max-max* model, and we set the evacuation time  $\tau$  to be sufficiently large so that there is no evacuation. We set the upper and lower thresholds to the same value,  $t_l(u) = t_h(u) = t(u)$  so that there is no Undecided state. We do not change any of the trust values along any of the edges, or between sources and nodes.

**Instance  $\mathbb{S}$** 

Graph  $G = (V, E)$  specifying the network for the diffusion.

Source information values and budgets,  $(I_1, B_1), \dots, (I_K, B_K)$

Trust values  $\alpha(u_i, u_j)$  for all edges  $(u_i, u_j) \in E$ .

Trust values  $\alpha(k, u_j)$  between each source  $k$  and each node  $u_j \in V$ .

Diffusion parameters  $\lambda_d = \lambda_s = 0, \tau \rightarrow \infty$

Lower and upper thresholds  $t_l(u) = t_h(u) = t(u)$  for each node  $u \in V$ .

**Desired output:** seed sets  $\psi_1, \dots, \psi_K$ .

Figure 5: The simplified instance  $\mathbb{S}$  for general instance  $\mathbb{G}$ .

When  $\lambda_s > 0$  or  $\lambda_d > 0$  the diffusion will generally be faster because taking the sum will tend to inflate information values at the nodes. To compensate for this, we need to raise the thresholds, to get a better approximation to  $\mathbb{G}$ . We treat  $t(u)$  as tunable parameters in the simplified model, and we may alter the values to get different instances of the simplified model. We will discuss how to choose these thresholds shortly. An instance of the simplified model is summarized in Figure 5. We distinguish quantities in the simplified model using a subscript  $\mathbb{S}$  from quantities in the general model with a subscript  $\mathbb{G}$ ; for example,  $\Gamma_{\mathbb{G}}(\psi_1, \dots, \psi_K)$  is the coverage function for a seeding in the general instance, and  $\Gamma_{\mathbb{S}}(\psi_1, \dots, \psi_K)$  is the coverage function for the same seeding in the simplified instance.

**Further Simplifying the Model** We can further simplify the model by insisting on a single source with information value  $I$  and budget equal to the total budget of the  $K$  sources,  $B_{\mathbb{S}} = B_1 + \dots + B_K$ . The information value  $I$  is chosen as the weighted information value of the  $K$  sources, where the weights are the number of seeds allocated to each source:

$$I_{\mathbb{S}} = \frac{1}{\sum_{k=1}^K B_k} \cdot \sum_{k=1}^K I_k \cdot B_k.$$

With  $K$  sources, every node  $u$  has a trust weight  $\alpha(k, u)$  for each source. These are combined into a single average trust weight with a single source:

$$\alpha_{\mathbb{S}}(1, u) = \frac{1}{K} \sum_{i=k}^K \alpha_{\mathbb{G}}(k, u).$$

All trust values between two nodes in the graph are unchanged. We will discuss next how to select the thresholds  $t_{\mathbb{S}}(u)$ . By considering different thresholds at the nodes, we are able to generate a variety of seedings, each of which is near optimal for a slightly different instance of the simplified model.

In this further simplified model, a single source will seed  $B_S$  nodes giving a seeding  $\psi_S$ . To convert it into a seeding for the general instance  $\mathbb{G}$  we need to assign each seeded node to one of the  $K$  sources in  $\mathbb{G}$ , so  $\psi_S \rightarrow \psi_1, \dots, \psi_K$ . In principle, one could try to optimize this assignment. For simplicity, we just pick an arbitrary assignment (for example a random assignment) of the seeded vertices to the  $K$  sources that respects the constraint  $|\psi_k| \leq B_k$ .

#### 4.5 Using Thresholds $t_S(u)$ to Generate Different Seedings

For different choices of  $t_S(u)$  with  $u \in V$ , we get different instances of the simplified model:

$$\{t_S(u)\} \xrightarrow{\text{greedy}} \psi_S \xrightarrow{\text{partition}} \psi_1, \dots, \psi_K \xrightarrow{\text{evaluate}} \Gamma_{\mathbb{G}}(\psi_1, \dots, \psi_K).$$

By evaluating according to  $\Gamma_{\mathbb{G}}$  in this way, we can use different settings of the thresholds in the simplified model to explore different candidate seedings for the general instance  $\mathbb{G}$  in a more intelligent way than the pure Brute Force approach.

There are several ways to choose the thresholds in the simplified model to get different seedings for the general instance. For simplicity and computational efficiency, we assume homogeneous thresholds, so every node has the same upper and lower threshold, and so  $t_S(u)$  becomes just  $t_S$ . We choose  $t_S$  from a set of thresholds  $\Omega = \{t_1, t_2, \dots, t_c\}$  where  $0 \leq t_i \leq 1$  to generate  $c$  simplified instances; this in turn will produce  $c$  seedings to be evaluated with  $\Gamma_{\mathbb{G}}$ . We choose  $\Omega$  and  $c$  as follows.

**Homogeneous Approximation to Trust in Instance  $\mathbb{G}$ .** To construct  $\Omega$ , we imagine all trust weights in instance  $\mathbb{G}$  are  $\alpha_{avg}$ , the average trust weight in the network for instance  $\mathbb{G}$  (this is just for the intuition on how we generate  $\Omega$ ). Let  $t_{min}^{\mathbb{G}}$  be the minimum lower threshold over all nodes in the general instance  $\mathbb{G}$ ; similarly  $t_{max}^{\mathbb{G}}$  is the maximum upper threshold. All the thresholds in  $\Omega$  will satisfy  $t_{min} \leq t \leq t_{max}$ . The highest threshold we need to consider is  $\alpha_{avg} \cdot I$  because no information value at a source can be higher than  $I$ , so if the threshold is higher, no node can convert to Believer. Every next hop reduces the information value by a factor of  $\alpha_{avg}$  so all thresholds which are in the interval  $[\alpha_{avg}^2 \cdot I, \alpha_{avg} \cdot I)$  are equivalent, and we only need consider one of them,  $\alpha_{avg}^2 \cdot I$ . We can continue this logic which implies that we only need to consider thresholds (in the simplified model) of the form  $t_i = \alpha_{avg}^i \cdot I$ . Further restricting the thresholds to be in the range  $[t_{min}, t_{max}]$  results in our choice of the threshold set  $\Omega$ :

$$\Omega = \{t_i \mid t_i = \alpha_{avg}^i \cdot I; t_{min} \leq t_i \leq t_{max}\} \cup \{t_{min}, t_{max}\}.$$

The number of thresholds  $c \approx 2 + O(\log t_{min} / \log \alpha_{avg})$ .

**Two-Level Approximation to Trust in Instance  $\mathbb{G}$ .** A better approximation to the trust values in  $\mathbb{G}$  can be obtained by clustering the trust weights into high values  $\alpha_{high}$  and low values  $\alpha_{low}$ . This leads to a

larger set of thresholds in  $\Omega$  for the simplified model. The general idea is the same. We choose thresholds as the possible information values at the nodes. This set of information values is

$$\{I; \alpha_{high} \cdot I, \alpha_{low} \cdot I; \alpha_{high}^2 \cdot I, \alpha_{high}\alpha_{low} \cdot I, \alpha_{low}^2 \cdot I; \dots\}$$

In general, we can construct the set  $\Omega$  iteratively as follows:

```

1:  $\Omega \leftarrow \{I\}$ .
2: while  $\alpha_{high} \cdot \max_{t \in \Omega} t \geq t_{low}$  do
3:    $\Omega_{add} \leftarrow \emptyset$ .
4:   for every  $t \in \Omega$  with  $\alpha_{high} \cdot t \geq t_{min}$  do
5:     add  $t$  to  $\Omega_{add}$ .
6:   for every  $t \in \Omega$  with  $\alpha_{low} \cdot t \geq t_{min}$  do
7:     add  $t$  to  $\Omega_{add}$ .
8:    $\Omega \leftarrow \Omega \cup \Omega_{add}$ 

```

This algorithm easily generalizes to approximating the trust in the instance  $\mathbb{G}$  by more than two trust levels. The size  $c$  of the resulting set of thresholds  $\Omega$  depends on how quickly information decays in the network; faster information decay leads to smaller  $c$ .

## 4.6 Running Time of Projected Greedy

The basic methodology of *Projected Greedy* includes constructing  $c$  instances of the simplified model and computing the near-optimal seeds in the simplified instance  $\mathbb{S}$ , converting that solution to a seeding for the general instance  $\mathbb{G}$  and finally evaluating  $\Gamma_{\mathbb{G}}$  for the seeding. The running time is therefore

$$c \cdot O(T_{greedy} + T_{\mathbb{G}}),$$

where  $T_{greedy}$  is the time to construct the solution for the simplified instance using the greedy algorithm. Our implementation of the greedy algorithm for the simplified instance exploits the monotone and submodular nature of  $\Gamma_{\mathbb{S}}$ . We create a special data structure  $\delta$  where  $\delta(u) \subseteq V$  is the set of nodes that convert  $u$  when they alone are added to the seed set. While calculating  $\delta$ , we also create array  $N$  where  $N(u)$  is the number of nodes that  $u$  (if added to the seed set) converts to the Believed state. For sparse graphs, building  $\delta$  and  $N$  takes  $O(n^2 \log n)$  time, where  $n = |V|$ . After selecting the node that gives the best improvement in  $\Gamma_{\mathbb{S}}$ , we only need to update the array  $N$  using  $\delta$ . This means we do not need to recalculate  $N$  at every step, and the update takes only  $O(n)$  time. The process of updating array  $N$  depends on properties of  $\Gamma_{\mathbb{S}}$ , such as monotonicity. Thus, for sparse graphs,

$$T_{greedy} = O(n^2 \log n + Bn).$$

When  $B = \Theta(n)$ , on sparse graphs, the total running time is typically  $c \cdot O(n^2 \log n)$ . To give an idea about how large the value of  $c$  is, consider an extreme case with homogenous trust values where  $t_{min} = 0.01$ ,  $t_{max} = 0.99$  and  $\alpha_{avg} = 0.9$ . In this extreme scenario  $c = 44$ , which gives a running time that is orders of magnitude better than *Actual Greedy* when  $n$  is in the millions.

Pre-calculating the data structure  $\delta$  can have worst case space complexity of  $O(n^2)$ . To avoid this, we use a hybrid strategy for computing seed sets in *Projected Greedy*. In this strategy  $\delta$  is not stored initially and the array  $N$  is recalculated each time the best node is to be selected. We use a technique in [29] according to which array  $N$  is stored as a priority queue and each selection requires recalculation of  $N(u)$  for only a small number of nodes  $u$ . Only after a certain threshold is breached do we switch to populating  $\delta$  for all nodes. Since nodes that convert the highest number of other nodes to Believed state have already been selected at this stage, it helps in reducing the size of  $\delta(u)$  for each of the remaining nodes. Not only does this help reduce memory usage but it also does not affect the running time adversely. Ultimately, *Projected Greedy* is more efficient by a factor of about  $n$  compared to *Actual Greedy*. In large graphs, this is a significant improvement and can sometimes be the difference between feasibility and intractability.

## 5 Experiment Design

In order to compare the Greedy heuristic with other seeding strategies, we simulate the spread of evacuation warnings in a social network. The simulation of diffusion is carried out on different types of networks with several parameter values.

### 5.1 Networks

We used three different network structures, each with 100,000 nodes with average degree approximately 4. Frequently, in real world scenarios individuals form groups based on race, ethnicity, nationality, etc. Not only are individuals within the same social group more likely to be acquainted with each other, they are also inclined to place more trust with people in the same group as theirs. Since these connections and trust disparities play an important role in the flow of information, we model social groups by dividing the population into 2 groups in the following networks.

**Scale-free Graphs:** We use the Albert-Barabasi model for generating random scale-free networks using the preferential attachment model [2, 5]. The graph thus produced follows a power law distribution with exponent approximately  $-2.9$ , that is  $P(x) \propto x^{-2.9}$  ( $P(x)$  is the fraction of nodes having degree  $x$ ). Once the graph is created, 50,000 nodes are randomly assigned to one group and the rest to the other. Scale-free graphs do not take into consideration the fact that nodes within a group are more likely to communicate

with each other.

**Random Group Model:** Nodes are randomly assigned into 2 groups of size 50,000 nodes. The probability that 2 nodes are connected (edge probability) depends on whether they belong to the same group or not. If two nodes belong to the same group then the edge probability is  $p_s$  while if they belong to different groups, the edge probability is  $p_d$ . Here  $p_s = 2 * p_d$  (more connections within a group as between groups). The probability  $p_d$  is chosen so that the average degree is 4.

**San Diego Network:** This is a random geometric graph that is constructed from actual demographic data in the San Diego area [20, 23]. Since there is a large population of Hispanics, we consider two groups: Hispanics and Non-hispanics. The population of nodes belonging to each group is based on the racial demographics of the region. Also, the edge probability between two nodes depends not only on the group they belong to, but also the physical distance between them. For example the edge probability between two nodes belonging to the same group is larger if they live close to each other than if they live far apart. Again the edge probabilities are chosen to have average degree 4. The details of this model are given in [20, 23].

## 5.2 Node Characteristics

Since evacuation is a high cost action, nodes would not evacuate without a significant amount of information. This can be modeled with high upper thresholds  $t_h(u)$ . Also, with such a high risk situation as evacuation, individuals may be proactive. That is to say, they may be willing to put more effort in collecting information. This can be modeled by decreasing the lower thresholds  $t_l(u)$ .

In our experiments, we simulate 3 different threshold value pairs for all nodes:  $(t_l = 0.2, t_h = 0.3)$ ,  $(t_l = 0.15, t_h = 0.55)$ , and  $(t_l = 0.4, t_h = 0.5)$ . We set the evacuation time to  $\tau = 5$  time-steps (so, all nodes leave the network 5 time steps after they are converted to their Believed state).

## 5.3 Trust Scenarios

Since nodes are split into 2 groups, there are 2 kinds of edges in the graph. The first type of edge is incident on nodes from the same group (denoted type  $A$  edges). The second type of edge is incident on nodes from different groups (denoted type  $B$  edges). Based on the trust values on these edges we have two trust scenarios. In each scenario, we set the average trust on the edges to be  $\alpha$ .

**Homogenous trust:** All edges have the same trust value. This models situations when no social groups exist. The trust value on every edge is  $\alpha$ .

**Group Variable trust:** The trust value on type  $A$  edges is  $\alpha + \varepsilon$  where  $\varepsilon > 0$ . The trust value on type  $B$  edges is chosen so that the average trust is  $\alpha$ . So the trust on type  $B$  edges will be less than type  $A$  edges, which models social groups that are more trusting of their own group than outsiders. We used  $\alpha = 0.7$  and  $\varepsilon = 0.05$  for our simulations.

## 5.4 Seeding Algorithms

We have 5 trustworthy sources each with information value  $I = 0.95$  and trust value  $\alpha(k, u) = 0.9$  for all  $u \in V$ . We look at scenarios in which between 5% – 50% of nodes are seeded in total, with each source seeding an equal number of nodes. We compare the following algorithms for generating the seeding of total size  $B$ .

**Random:** Randomly select  $B$  nodes and arbitrarily assign these nodes to the  $K$  sources.

**High Degree:** Select the  $B$  highest degree nodes and arbitrarily assign them to the  $K$  sources. Here, the degree for node  $u$  is the total outgoing trust weight:

$$degree(u) = \sum_{(u, u') \in E} \alpha(u, u')$$

**Projected Greedy heuristic:** Seeds are generated according to the Projected Greedy heuristic described in Section 4.

## 5.5 Parameters

For all our experiments we use  $\lambda_s = 0$ . This means nodes always chose the maximum value while combining information from the same source. We ran simulations for different choices of  $\lambda_d \in \{0.0, 0.05, 0.1, 0.2\}$ .

Lastly, in real life scenarios, communication between nodes is not likely to succeed every time. Hence, when a node that is not a source, queries or propagates an information set to another node, it will succeed with some probability  $p$ ; we set  $p = 0.75$  in our experiments.

# 6 Results and Discussion

We ran each simulation for 50 steps and repeat it 100 times. We observed that 50 steps are enough for the diffusion process to conclude. Since the networks are generated randomly, we repeat the whole simulation on at least 10 instances of the graph. We use the average number of nodes evacuated as our measure of performance. The standard deviation (error bar) due to the randomness is extremely small; in fact the error bars are not even visible in our plots.



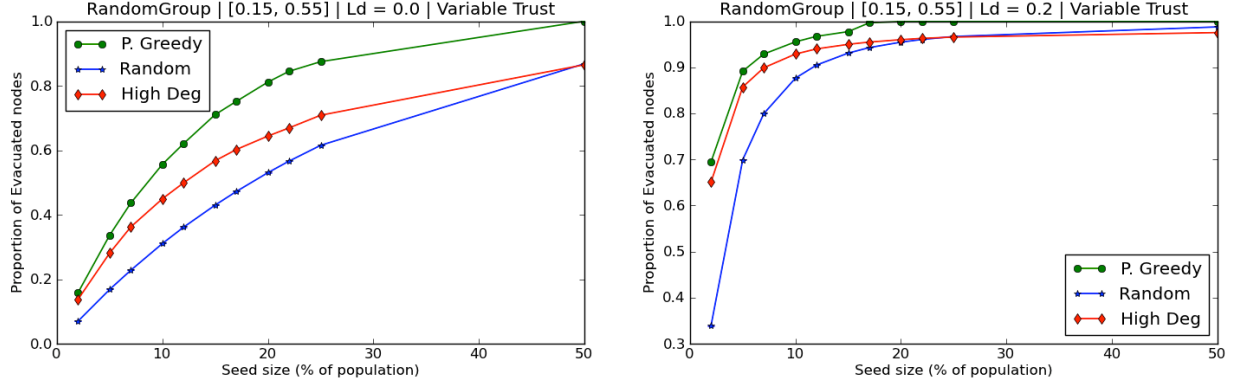


Figure 6: Random Group network with group variable trust (high trust edges within a group). Thresholds are  $t_l = 0.15$ ,  $t_h = 0.55$ ; We show  $\lambda_d = 0.0$  (left) versus  $\lambda_d = 0.2$  (right) for various total seed budget. Note that for plot on the right, the  $y$ -axis range is  $[0.0, 1.0]$  while for the right plot it is  $[0.3, 1.0]$ . The Projected Greedy dominates the other algorithms.

## 6.1 Seed Size

We start with the performance of the three seeding strategies (random, high-degree and Projected Greedy) as the total seeding budget increases. Figure 6 shows the random group model and Figure 7 the San Diego network. Though we ran several different scenarios, we pick these two as both relevant and representative of the typical nature of the results.

For both the random graph model and the San Diego network, the Projected Greedy heuristic performs consistently better than both Random and High Degree seeding strategies. When  $\lambda_d = 0.0$ , Projected Greedy and High Degree are comparable for a small seed budget, but as the seeding budget increases, so does the gap between their performance. This is likely because there are two competing effects: one should seed influential nodes, and one should spread out the seeds (i.e., not have the seed nodes be too close to each other, so that they will convert more nodes in total). When you have a few seeds, influence is more important, and Projected Greedy and high-degree are achieving that goal comparably. When you have many seeds, it now becomes more important to spread the seeds out, and high-degree ignores how spread out the seeds are, whereas Projected Greedy takes that into account. This also explains why, eventually, even random seeding becomes better than high degree: with more seeds, the need to spread out the seeds wins. This fact is even more pronounced when  $\lambda_d = 0.2$  which makes sense because many spread out seeds can have even more impact due to the summing effect present when  $\lambda_d > 0$ . This effect is even stronger in the San Diego network (Figure 7), probably due to the structure of San Diego network. Such typical social networks have a number of small dense clusters that weakly connect to each other. This makes selecting

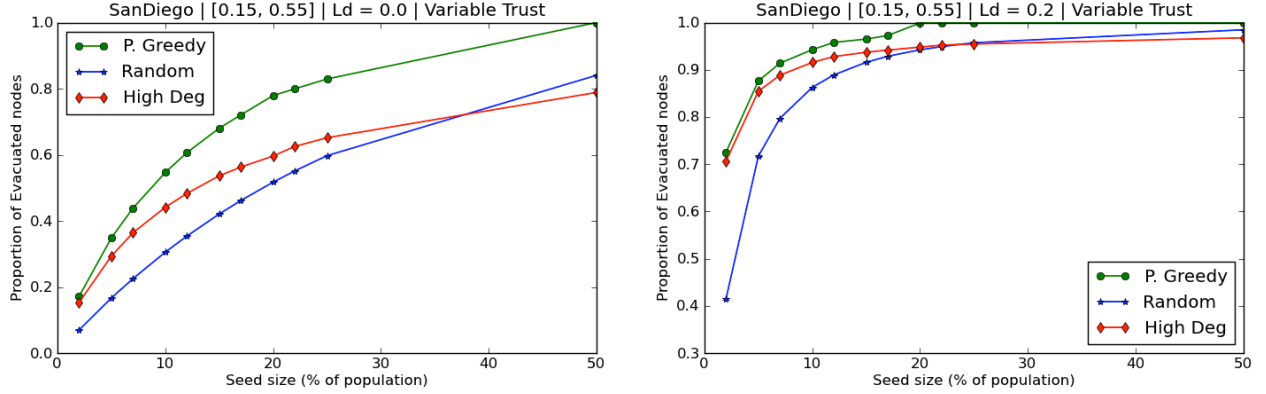


Figure 7: San Diego network, with group variable trust (high trust edges within a group). Thresholds are  $t_l = 0.15$ ,  $t_h = 0.55$ ; We show  $\lambda_d = 0.0$  (left) versus  $\lambda_d = 0.2$  (right) for various total seed budget. Note that for plot on the right, the  $y$ -axis range is  $[0.3, 1.0]$  while for the right plot it is  $[0.0, 1.0]$ . Projected Greedy clearly dominates the other seeding strategies.

high degree nodes from the same cluster a bad idea.

The comparative advantage of Projected Greedy is higher with  $\lambda_d = 0$  for two reasons. The first is that the simplified instance  $\mathbb{S}$  on which the seed set is near optimal is a better approximation to the true instance  $\mathbb{G}$  when  $\lambda_d = 0$ . Second, when you have some element of summing in the diffusion process, all algorithms will improve, and hence their differences will diminish. Nevertheless, Projected Greedy still outperforms even for  $\lambda_d = 0.2$ , managing to convert almost every node to believed state with seed set size as low as 20%; even at 50% seeds, the other algorithms cannot achieve this.

## 6.2 Cross section of Results

We give a comprehensive cross section of the results in Table 1 for all three types of networks, both homogeneous and group-variable trust, and for a variety of trust thresholds. We quantify the performance in terms of the regret: for a given scenario, the best performing algorithm has regret zero, and otherwise,

$$regret = \frac{\Gamma_{best} - \Gamma}{\Gamma_{best}} \times 100\%,$$

where  $\Gamma_{best}$  is the number of nodes evacuated by the best algorithm for that scenario, and  $\Gamma$  is the number of nodes evacuated by the algorithm whose regret is being computed. A good algorithm would always have low regret. All results are for a total seeding size of 5%. The actual fraction of the network evacuated is shown in Table 2 for a particular instance.

In general, we make the following observations. When the number of evacuated nodes is large, which occurs either with low thresholds or a high value for  $\lambda_d$  (more information aggregation), the performance edge

			Homogenous Trust			Variable Trust		
			R	HD	PG	R	HD	PG
SF	$t_l = 0.2, t_h = 0.3$	$\lambda_d : 0.0$	25.76	1.33	<b>0.00</b>	4.32	0.08	<b>0.00</b>
		$\lambda_d : 0.2$	<b>0.00</b>	<b>0.00</b>	<b>0.00</b>	<b>0.00</b>	<b>0.00</b>	<b>0.00</b>
	$t_l = 0.15, t_h = 0.55$	$\lambda_d : 0.0$	70.76	3.23	<b>0.00</b>	71.00	3.20	<b>0.00</b>
		$\lambda_d : 0.2$	25.75	1.33	<b>0.00</b>	14.39	0.63	<b>0.00</b>
RG	$t_l = 0.2, t_h = 0.3$	$\lambda_d : 0.0$	31.85	11.16	<b>0.00</b>	16.21	3.87	<b>0.00</b>
		$\lambda_d : 0.2$	2.40	2.03	<b>0.00</b>	2.77	2.30	<b>0.00</b>
	$t_l = 0.15, t_h = 0.55$	$\lambda_d : 0.0$	46.46	10.94	<b>0.00</b>	50.06	16.48	<b>0.00</b>
		$\lambda_d : 0.2$	32.20	10.64	<b>0.00</b>	21.68	3.87	<b>0.00</b>
SD	$t_l = 0.2, t_h = 0.3$	$\lambda_d : 0.0$	31.85	11.02	<b>0.00</b>	14.09	2.91	<b>0.00</b>
		$\lambda_d : 0.2$	2.68	2.33	<b>0.00</b>	3.32	2.76	<b>0.00</b>
	$t_l = 0.15, t_h = 0.55$	$\lambda_d : 0.0$	47.28	11.39	<b>0.00</b>	52.14	15.60	<b>0.00</b>
		$\lambda_d : 0.2$	32.19	10.52	<b>0.00</b>	18.16	2.41	<b>0.00</b>

Table 1: Cross section of results. The table shows the relative regret of an algorithm for a particular scenario, a regret of **0** indicating the best performing algorithm. The Project Greedy heuristic significantly dominates all the other algorithms in almost every scenario. SF=scale free network; RG=random group model; SD=San Diego network.

delivered by Projected Greedy is smallest. In fact when  $\lambda_d$  increases beyond 0.2, the effect of information aggregation quickly takes over and there is little difference between the performance of the three algorithms as almost all nodes are evacuated with ease. When the thresholds are high or  $\lambda_d \approx 0$  the performance edge delivered by Projected Greedy is quite significant, in accordance to results shown in Figures 6 and 7. Projected Greedy's advantage increases with more seeds.

Note that in scale free graphs, high-degree is comparable to Projected Greedy but random significantly underperforms. This is because in such networks there are a few extremely important nodes, and it is essential to include these nodes in the seed set, which is unlikely to happen with random seeding.

		Homogenous Trust			Variable Trust		
		R	HD	PG	R	HD	PG
$t_l = 0.15, t_h = 0.55 \lambda_d : 0.0$	SF	20.82	68.89	71.19	20.65	68.92	71.20
	RG	22.26	37.03	41.58	16.85	28.17	33.73
	SD	22.19	37.30	42.09	16.72	29.48	34.93

Table 2: Fraction of the network evacuated for a particular scenario. Random has roughly the same performance on all networks and is a significant under performer on scale free networks.

### 6.3 Threshold Selection

As described in Section 4, in order to select the best seed set we carry out simulations over a range of values for threshold  $t$  in the simplified model. Specifically we repeat the simulation  $c$  times. It is interesting to see how the coverage  $\Gamma_{\mathbb{G}}$  changes as we change the thresholds in the simplified model  $\mathbb{S}$ . In particular, how the best threshold  $t_{opt}$  in the simplified model (i.e. the threshold that gives the closest approximation to  $\mathbb{G}$ ) depends on the parameters of the general model, in particular  $(\lambda_s, \lambda_d)$ , the lower and upper node thresholds  $(t_l, t_u)$  and transmission probability  $p$ .

Our simulation results show that there is an interesting relationship between  $t_{opt}$  (the best threshold to choose in  $\mathbb{S}$  and  $t_u$  (the upper thresholds in  $\mathbb{G}$ ) under different values of  $\lambda_d$ . In order to observe this relationship, we perform simulations on a generalized version of the San Diego network. In this generalized version, edge trust values and node thresholds are selected uniformly at random from a range instead of being fixed to specific values. Thus the generalized network tries to incorporate variations observed in real networks. The network used for simulation has the following parameter values. For edges between nodes belonging to the same group, trust values are selected uniformly at random from the range  $[0.7, 0.8]$ . For edges between nodes belonging to different groups, trust values are selected uniformly at random from the range  $[\alpha_{low} - 0.5, \alpha_{low} + 0.5]$ . Here  $\alpha_{low}$  is selected such that the expected value of edge trusts in the network is 0.7. Similarly, for every  $u \in V$ ,  $t_l(u)$  is selected uniformly at random from range  $[0.1, 0.2]$  and  $t_u(u)$  from range  $[0.5, 0.6]$ . For the probability of transmission  $p = 0.75$  and seed set size 5% of total nodes, Figure 8 shows the proportion of nodes evacuated with threshold values for Projected greedy in range  $[0.1, 0.6]$  and  $\lambda_d = [0.0, 0.05, 0.1, 0.15, 0.2]$ .

It is intuitive to expect that  $t_{opt}$  will be close to  $E[t_u]$  since nodes get converted to Believed state only after information value crosses  $t_u$ . This is exactly what we see in the case where  $\lambda_d = 0.0$ . The optimal threshold  $t_{opt}$  is very close to  $E[t_u] = 0.55$ . But as  $\lambda_d$  increases we see a gradual decrease in the value of  $t_{opt}$ . As  $\lambda_d$  increases, information gets aggregated and the fused information value becomes progressively

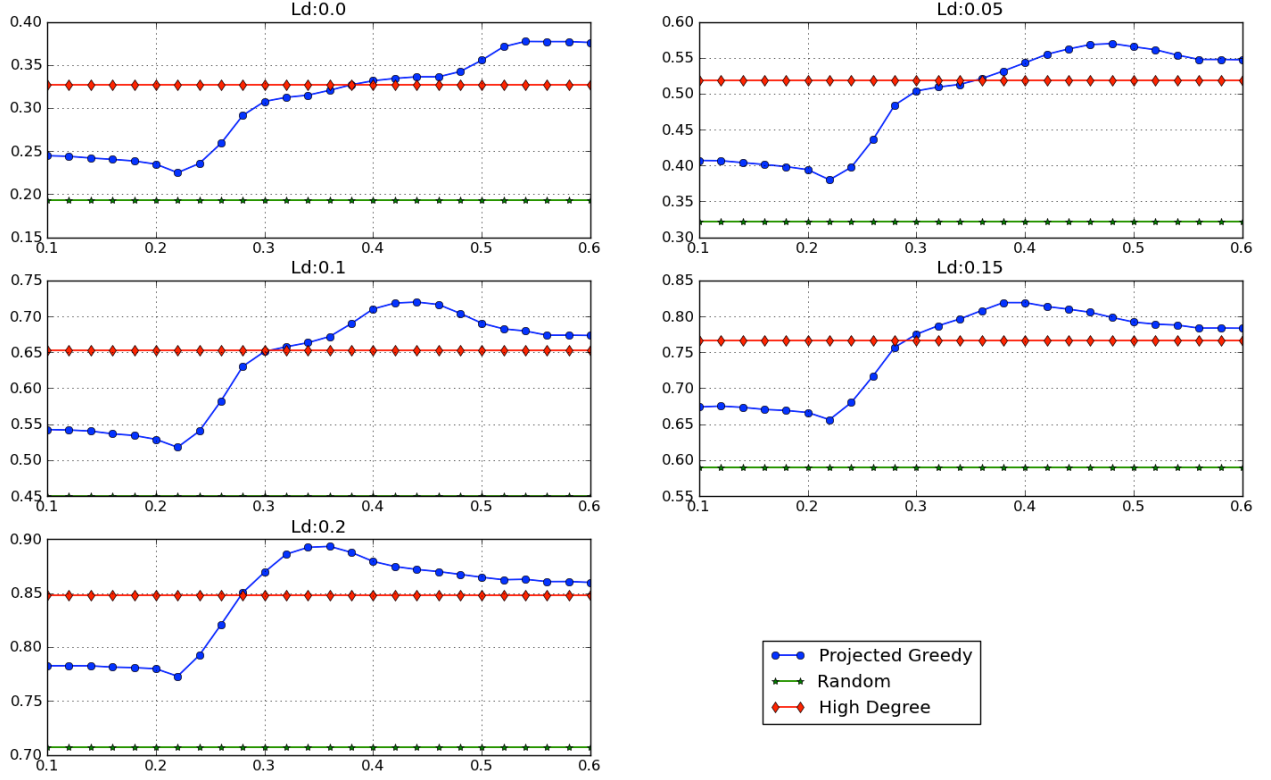


Figure 8: Plots showing change in the coverage in the general instance  $\mathbb{G}$  obtained from using different node thresholds in the simplified instance  $\mathbb{S}$ . We are interested in how the best threshold in the simplified instance  $\mathbb{S}$  changes (the threshold that gives the maximum coverage in  $\mathbb{G}$ ) as we increase  $\lambda_d$ . The  $x$ -axis shows the threshold in the simplified instance  $\mathbb{S}$  and the  $y$ -axis shows the proportion of nodes evacuated in the general input instance  $\mathbb{G}$ . The plots are labeled with the value of  $\lambda_d$ .

more successful at breaching  $t_u$ . In other words, the general model starts behaving like a simple model with a smaller  $t_u$  value where, for the same amount of initial information, it is comparatively easier to convert nodes. Table 3 shows how the value of  $t_{opt}$  reduces with increasing  $\lambda_d$ . The rate at which  $t_{opt}$  decreases may depend upon factors such as the type of network (i.e. network structure),  $t_l$  etc. These are interesting questions for future work.

## References

- [1] B.E. Aguirre. *Planning, warning, evacuation, and search and rescue: a review of the social science research literature*. Hazard Reduction & Recovery Center, College Station, TX, 1994.

Table 3:  $t_{opt}$  decreases as  $\lambda_d$  increases.

$\lambda_d$	$t_{opt}$
0.0	0.54
0.05	0.48
0.1	0.44
0.15	0.38
0.2	0.36

- [2] R. Albert and A.-L. Barabási. Statistical mechanics of complex networks. *Rev. of Modern Phys.*, 74:47–97, Jan 2002.
- [3] H.M. Amman, L. Tesfatsion, K.L. Judd, D.A. Kendrick, and J. Rust. *Handbook of Computational Economics: Agent-Based Computational Economics*, chapter Agent-based models of innovation and technological change. Handbooks in Economics. Elsevier, 2006.
- [4] Norman Bailey. *The Mathematical Theory of Infectious Diseases and its Applications*. Griffin, London, 1975.
- [5] A. L. Barabási and R. Albert. Emergence of scaling in random networks. *Science*, 286:509–512, 1999.
- [6] Frank M. Bass. A new product growth for model consumer durables. *Management Science*, 15(5):215–227, January 1969.
- [7] Thomas Berger. Agent-based spatial models applied to agriculture: a simulation tool for technology diffusion, resource use changes and policy analysis. *Agricultural Economics*, 25(23):245 – 260, 2001. Increasing Efficiency in Production, Research, Markets and Environmental Management. Selected and edited papers presented during the XXIV Conference of the International Association of Agricultural Economists.
- [8] E. Bonabeau. Agent-based modeling: Methods and techniques for simulating human systems. *Proc. National Academy of Sciences*, 99(Suppl 3):7280, 2002.
- [9] Gruia Calinescu, Chandra Chekuri, Martin Pál, and Jan Vondrák. Maximizing a monotone submodular function subject to a matroid constraint. *SIAM J. on Comput.*, 40(6):1740–1766, 2011.

- [10] R Chellappa and A Jain, editors. *Markov Random Fields: Theory and Application*. Academic Press, Waltham, Mass, 1993.
- [11] K. Chopra and W. A. Wallace. Modeling relationships among multiple graphical structures. *Computational and Mathematical Organization Theory*, 6(4):361–380, December 2000.
- [12] K. Chopra and W. A. Wallace. Trust in electronic environments. In *Proc. 36th Hawaii International Conference on Systems Science*, Hawaii, 2003.
- [13] Sebastiano A. Delre, Wander Jager, and Marco A. Janssen. Diffusion dynamics in small-world networks with heterogeneous consumers. *Comput. Math. Organ. Theory*, 13(2):185–202, June 2007.
- [14] Amy Wenxuan Ding. A theoretical model of public response to the homeland security advisory system. *The J. of Defense Modeling and Simulation: Applicat., Methodology, Technology*, 3(1):45–55, January 2006.
- [15] M. L. Fisher, G. L. Nemhauser, and L. A. Wolsey. An analysis of approximations for maximizing submodular set functions, II. In M. L. Balinski and A. J. Hoffman, editors, *Polyhedral Combinatorics*, volume 8 of *Mathematical Programming Studies*, pages 73–87. Springer Berlin Heidelberg, Philadelphia, PA, 1978.
- [16] Joseph H. Golden and Christopher R. Adams. The tornado problem: Forecast, warning, and response. *Natural Hazards Review*, 1(2):107–118, 2000.
- [17] Mark Granovetter. Threshold models of collective behavior. *American Journal of Sociology*, 83(6):pp. 1420–1443, 1978.
- [18] Daniel Gruhl, R. Guha, David Liben-Nowell, and Andrew Tomkins. Information diffusion through blogspace. In *Proc. of the 13th Int. Conf. on World Wide Web*, WWW '04, pages 491–501, New York, NY, 2004. ACM.
- [19] Shawndra Hill, Foster Provost, and Chris Volinsky. Network-based marketing: Identifying likely adopters via consumer networks. *Stat. Sci.*, 21:256–276, 2006.
- [20] Cindy Hui. *Modeling the diffusion of information in dynamic social and communication networks*. PhD thesis, Decision Sci. and Eng. Syst., Rensselaer Poly. Inst., Troy, NY, 2011.
- [21] Cindy Hui, Mark Goldberg, Malik Magdon-Ismail, and William A. Wallace. Agent-based simulation of the diffusion of warnings. In *Proceedings of the 2010 Spring Simulation Multiconference*, SpringSim '10, pages 9:1–9:8, San Diego, CA, USA, 2010. Society for Computer Simulation International.

- [22] Cindy Hui, Mark Goldberg, Malik Magdon-Ismail, and William A. Wallace. Simulating the diffusion of information: An agent-based modeling approach. *International Journal of Agent Technologies and Systems (IJATS)*, 2(3):31–46, 2010.
- [23] Cindy Hui, Malik Magdon-Ismail, William A. Wallace, and Mark Goldberg. Importance of ties in information diffusion. Poster presentation at Workshop on Information In Networks (WIN), Sept. 2010.
- [24] Jose Luis Iribarren and Esteban Moro. Information diffusion epidemics in social networks. *Phys. Rev. Lett.*, 103:038702, 2009.
- [25] Akshay Java, Pranam Kolari, Tim Finin, and Tim Oates. Modeling the Spread of Influence on the Blogosphere. Technical report, University of Maryland, Baltimore County, Baltimore, MD, March 2006.
- [26] Kari Kelton, Ken R. Fleischmann, and William A. Wallace. Trust in digital information. *J. Amer. Society for Information Science and Technology*, 57(3):363–374, 2008.
- [27] David Kempe, Jon Kleinberg, and Éva Tardos. Maximizing the spread of influence through a social network. In *Proc. of the 9th ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining*, KDD ’03, pages 137–146, New York, NY, 2003. ACM.
- [28] David Kempe, Jon Kleinberg, and Éva Tardos. Influential nodes in a diffusion model for social networks. *Automata Languages and Programming*, 3580(1):99–120, 2006.
- [29] Andreas Krause, Jure Leskovec, Carlos Guestrin, Jeanne VanBriesen, and Christos Faloutsos. Efficient sensor placement optimization for securing large water distribution networks. *Journal of Water Resources Planning and Management*, 134(6):516–526, November 2008. (Draft; full version available here).
- [30] Jure Leskovec, Lada A. Adamic, and Bernardo A. Huberman. The dynamics of viral marketing. *ACM Trans. Web*, 1, May 2007.
- [31] Jure Leskovec, Ajit Singh, and Jon Kleinberg. Patterns of influence in a recommendation network. In *Proc. of the 10th Pacific-Asia Conf. on Advances in Knowledge Discovery and Data Mining*, PAKDD’06, pages 380–389, Singapore, 2006. Springer-Verlag.
- [32] M. K. Lindell and C. S. Prater. Estimating evacuation time components: Lessons from nuclear power plants, hurricanes, and the first world trade center bombing. In *National Institute for Standards and Technology Workshop on Building Occupant Movement During Fire Emergencies*, pages 91–95, Gaithersburg, MD, 2004.



- [33] M. K. Lindell and Carla S. Prater. Critical behavioral assumptions in evacuation time estimate analysis for private vehicles: Examples from hurricane research and planning. *J. of Urban Planning and Develop.*, 133(1):18–29, 2007.
- [34] Michael K. Lindell, Carla S. Prater, and Walter Gillis Peacock. Organizational communication and decision making for hurricane emergencies. *Natural Hazards Rev.*, 8(3):50–60, 2007.
- [35] M.K. Lindell and R.W. Perry. *Communicating environmental risk in multiethnic communities*. Communicating effectively in multicultural contexts. Sage Publications, CA, 2004.
- [36] Tiejun Ma and Yoshiteru Nakamori. Agent-based modeling on technological innovation as an evolutionary process. *European Journal of Operational Research*, 166(3):741 – 755, 2005. *Advances in Complex Systems Modeling*.
- [37] Dennis Mileti. *Disasters by Design: A Reassessment of Natural Hazards in the United States*. The National Academies Press, Washington, DC, 1999.
- [38] P. Monge and N. Contractor. *Theories of Communication Networks*. Oxford University Press, New York, 2003.
- [39] F. Neri. Agent based simulation of information diffusion in a virtual market place. In *Intelligent Agent Technology, 2004. (IAT 2004). Proceedings. IEEE/WIC/ACM International Conference on*, pages 333 – 336, sept. 2004.
- [40] M. E. J. Newman. The spread of epidemic disease on networks. *Phys. Rev. Lett.*, 66:016128, 2002.
- [41] M A Nowak and R M May. *Virus Dynamics: Mathematical Principles of Immunology and Virology*, volume 291. Oxford University Press, Oxford, 2000.
- [42] Committee on Disaster Research in the Social Sciences: Future Challenges and National Research Council Opportunities. *Facing Hazards and Disasters: Understanding Human Dimensions*. The National Academies Press, Washinton, D.C., 2006.
- [43] L. Perez and S. Dragicevic. An agent-based approach for modeling dynamics of contagious disease spread. *International Journal of Health Geographics*, 8:50, 2009.
- [44] Ronald W. Perry and Michael K. Lindell. Understanding citizen response to disasters with implications for terrorism. *J. of Contingencies and Crisis Management*, 11(2):49–60, 2003.
- [45] Ronald W. Perry, Michael K. Lindell, and Marjorie R. Greene. Crisis communications: Ethnic differentials in interpreting and acting on disaster warnings. *Social Behavior and Personality: an Int. J.*, 10(1):97–104, 1982.

- [46] R.W. Perry and A.H. Mushkatel. *Minority Citizens in Disasters*. University of Georgia Press, Athens, GA, 2008.
- [47] G. O. Rogers and J. H. Sorensen. *Risk Analysis Prospects and Opportunities*, chapter Diffusion of Emergency Warning: Comparing Empirical and Simulation Results, pages 117 – 134. Plenum Press, New York, 1991.
- [48] Nina Schwarz and Andreas Ernst. Agent-based modeling of the diffusion of environmental innovations – an empirical approach. *Technological Forecasting and Social Change*, 76(4):497 – 511, 2009.
- [49] J H Sorensen. Hazard warning systems: review of 20 years of progress. *Natural Hazards Rev.*, 1(2):119–125, 2000.
- [50] J. H. Sorensen. Modeling human response to a chemical weapons accident. In *DIMACS Working Group on Modeling Social Responses to Bio-terrorism Involving Infectious Agents*, New Brunswick, NJ, 2003.
- [51] J. H. Sorensen and D. S. Mileti. Decision-making uncertainties in emergency warning system organizations. *Int. J. of Mass Emergencies and Disasters*, 5(1):33–61, 1987.
- [52] J.H. Sorensen, United States. Dept. of Energy, Martin Marietta Energy Systems Inc., and Oak Ridge National Laboratory. *Assessment of the need for dual indoor/outdoor warning systems and enhanced tone alert technologies in the Chemical Stockpile Emergency Preparedness Program*. Oak Ridge National Laboratory, Oak Ridge, Tennessee, 1992.
- [53] J.H. Sorensen, B.M. Vogt, and D.S. Mileti. *Evacuation: an assessment of planning and research*. Oak Ridge National Laboratory, Oak Ridge, Tennessee, 1987.
- [54] David Strang and Sarah A Soule. Diffusion in organizations and social movements: From hybrid corn to poison pills. *Annu. Review of Sociol.*, 24(1):265–290, 1998.
- [55] Thomas W. Valente. Network models of the diffusion of innovations. *Comput. & Math. Organization Theory*, 2:163–164, 1996.
- [56] Xiaojun Wan and Jianwu Yang. Learning information diffusion process on the web. In *Proc. of the 16th Int. Conf. on World Wide Web*, WWW '07, pages 1173–1174, New York, NY, 2007. ACM.
- [57] H Peyton Young. The diffusion of innovations in social networks. Economics Working Paper Archive 437, The Johns Hopkins University, Department of Economics, May 2000.